The Web IN4325 - Information Retrieval



Today

- The Web
- HITS & PageRank
- Sentiment Analysis
 - Naïve Bayes classifier
- Spam
 - Decision tree classifier
- The Social Web
 - Personalized search

next Monday Prof. Arjen de Vries on entities

next Wednesday

search engine clicks (ranking & advertising)



The Web graph





- The World Wide Web was formed in the early 1990's
 - Creator: Tim Berners-Lee
 - make documents easily available to anyone on the Internet (Web pages)
 - Easy access to such Web pages using a *browser*
- Early Web years
 - Full-text search engines (Altavista, Excite and Infoseek) vs.
 - Taxonomies populated with pages in categories (**ODP**, Yahoo! Directory)



The Web

 Today the size of the Web makes it impossible to find anything without the help of search engines

- Estimating the size of the Web is a research area by itself (by comparing the overlap between Web search engines output)
- In 1998 (Google): 24 million pages (at that time: academic engine)
- In 1999 (Altavista): >200 million pages
- In 2005: *indexed* Web estimated to be at 11.5 billion pages [7]
- Users view the Web through the lense of the search engine
- Pages not indexed (or ranked at low positions) by search engines are unlikely to be found by users





Research questions

Why do people link? How connected is the web graph? What properties does the graph have? How can we exploit those properties?

Web graph

- nodes: web pages, web sites, domains (dependent on abstraction level
- directed edges: hyperlinks

Graph structure in the Web

Broder et al., 1999 [4]

>2000 citations

- Insights important for
 - The design of crawling strategies
 - Understanding the sociology of content creation on the Web
 - Analyzing the behaviour of algorithms that rely on link information (e.g. HITS, PageRank)
 - Predicting the evolution of web structures
 - Predicting the emergence of new phenomena in the Web graph
- Data: Altavista crawl from 1999 with 200M pages and 1.5 billion links
- In/out-degrees follow power laws
 - Probability that a node has in/out-degree i is proportional to $\frac{1}{x}$, x > 1

Graph theory: connected components

A detour

• Strongly connected component (SCC): directed graph with a path from each node to every other node $\begin{pmatrix} A & B \\ A & 0 \end{pmatrix}$

Weakly connected component (WCC): directed graph with a path in the underlying undirected graph from each node to every other node
 (A)

2 weakly connected components

Graph theory: connected components

A detour

• Strongly connected component (SCC): directed graph with a path from each node to every other node $\begin{pmatrix} A \\ A \end{pmatrix}$

G = (V, E) $V = \{A, B, C, D\}$ $E = \{(A, D), (B, C), (C, A), (C, B), (C, D), (D, B)\}$ d(A, B) = 2, d(C, B) = 1, d(A, C) = 3

Weakly connected component (WCC): directed graph with a path in the underlying undirected graph from each node to every other node
 (A)

G = (V, E) $V = \{A, B, C, D\}$ $E = \{\{A, C\}, \{A, D\}, \{B, C\}, \{B, D\}, \{C, D\}\}$ d(A, B) = 2, d(C, B) = 1, d(A, C) = 1

² weakly connected components

Graph theory: diameter

A detour

• Graph diameter: longest shortest path in the graph

12

Graph structure in the Web

Broder et al., 1999 [4]

 "In a sense the web is much like a complicated organism, in which the local structure at a microscopic scale looks very regular like a biological cell, but the global structure exhibits interesting morphological structure (body and limbs) that are not obviously evident in the local structure."

Evolution of the Web

Fetterly et al., 2003 [8]

- How fast does the Web change?
 - Important to build effective crawlers (what portions of the index should be updated?
- Setup: weekly crawl of 150M pages once a week for 11 weeks (2002-2003)

Figure 1. Number of successful, completed, and attempted downloads per URL.

Source: [7]

ŤUDelft

Kleinberg, 1998 [12]

- Hyperlink-Induced Topic Search
- Intuition: two broad types of useful pages for ad hoc search queries
 - **Authoritative** pages: pages containing a lot of relevant information about the search topic
 - Wikipedia pages, etc.
 - High weight *a*(*p*)
 - **Hub** pages: pages containing a large number of useful hyperlinks pointing to pages with relevant content
 - Yahoo! Portal, Open Directory Project, etc.
 - High weight *h(p)*

- Root set (RS): retrieve the top 200 results for a given keyword query
- ② Base set (BS): expand RS by including all* pages that link to pages in RS or are linked-to by pages in RS
- Clean hyperlink structure (link removal between pages belonging to the same web site)
- Initialize all hub/authority weights to 1

 $q \rightarrow p$

5 Iteratively update hub/authority weights (& normalize)

$$a(p) = \sum h(q) \qquad h(p) = \sum a(q)$$

authority weight increased if good hubs point to \ensuremath{p}

hub weight of p increased if it points to good authorities

 $p \rightarrow q$

$$a(p) = \sum_{q \to p} h(q)$$

initialize :
$$\vec{a} = (a(1),...,a(5)) = (1,1,1,1,1)$$

first round :

$$a(1) = h(3) = 1 \propto \frac{1}{6}$$

$$a(2) = h(1) + h(3) + h(4) = 1 + 1 + 1 \propto$$

$$a(3) = h(5) = 1 \propto \frac{1}{6}$$

$$a(4) = h(3) = 1 \propto \frac{1}{6}$$

$$a(5) = 0$$

 $\frac{1}{2}$

$$h(p) = \sum_{p \to q} a(q)$$

initialize:

$$\vec{h} = (h(1), ..., h(5)) = (1, 1, 1, 1, 1)$$

first round:
 $h(1) = a(2) = 1 \propto \frac{1}{6}$
 $h(2) = 0$
 $h(3) = a(1) + a(2) + a(4) = 3 \propto \frac{1}{2}$
 $h(4) = a(2) = 1 \propto \frac{1}{6}$
 $h(5) = a(3) = 1 \propto \frac{1}{6}$

normalization after each round (sum to 1)

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 3 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$A \qquad \qquad \times \vec{a}_{init} = \vec{h}_{1}$$

$$\vec{h}_m = A\vec{a}_{m-1}$$
$$\vec{a}_m = A^T\vec{h}_{m-1}$$

$$\vec{h}_m = (A \times A^T) \vec{h}_{m-2}$$
$$\vec{a}_m = (A^T \times A) \vec{a}_{m-2}$$

Our toy example: weights across 8 iterations

HITS output: a list of top-scoring authority/hub pages for the given query

PageRank

Page et al., 1998 [5]

>4700 citations

- A **topic independent** approach to page importance
 - Computed once per crawl
- Every document of the corpus is assigned an importance score
 - In search: re-rank (or filter) results with a low PageRank score
- Simple idea: number of in-links importance
 - Page p₁ has 10 in-links and one of those is from yahoo.com, page p₂ has 50 in-links from obscure pages
- PageRank takes the importance of the page where the link originates into account

"To test the utility of PageRank for search, we built a web search engine called Google."

PageRank

Page et al., 1998 [5]

- Idea: if page p_x links to page p_y , then the creator of p_x implicitly transfers some importance to page p_y
 - yahoo.com is an important page, many pages point to it
 - pages linked to from yahoo.com are also likely to be important
- Each page distributes "importance" through its outlinks
- Simple PageRank (iteratively):

$$PageRank_{i+1}(v) = \sum_{u \to v} \frac{PageRank_i(u)}{N_u}$$
 particular page.
all nodes linking to v out-degree of node u

A page with many out-links has little influence on one

PageRank

Simplified formula

initialize PageRank vector \vec{R} $\vec{R} = (R(1), \dots, R(4)) = (0.25, 0.25, 0.25, 0.25)$ 0.33 $W^1 \times \vec{R}' = \begin{vmatrix} 0.46 \\ 0.13 \\ 0.02 \end{vmatrix}$ 0.08 $W^{16} \times \vec{R}' = \begin{pmatrix} 0.40 \\ 0.33 \\ 0.20 \\ 0 \end{pmatrix}$ 0.50 $W^2 \times \vec{R}' = \begin{vmatrix} 0.29 \\ 0.17 \\ 0.17 \end{vmatrix}$ PageRank vector 0.35 0.40converges $W^{17} \times \vec{R}' =$ $W^{3} \times \vec{R}' = \begin{vmatrix} 0.35 \\ 0.25 \end{vmatrix}$ 0.34 eventually 0.20 0.06 0.07

 $PageRank_i = W \times PageRank_{i-1}$

PageRank

Reality

 $\begin{bmatrix} 5 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 3 \end{bmatrix} W = \left[\begin{array}{cccc} 0 & 0 & 0 & 0 \\ 1 & 0 & 1/2 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1/2 & 0 \end{array} \right]$

PageRank applications

Search

- Re-rank the top retrieved documents of a content retrieval technique according to the pages' PageRank score
- Filter out pages with low PageRank scores
- Personalized PageRank
 - Instead of random teleporting, bias the teleport locations
- PageRank as future inlink count predictor
 - Re-order crawling list accordingly (crawl better pages first)

Opinion mining

Opinion mining

Top Critic

- Opinion mining is the automatic extraction of subjective information from documents
- Applications
 - Automatic review aggregation websites (what do people think about X?)
 - Recognition of user rating errors: review vs. rating polarity
 - Flame detection & detection of sensible content (e.g. accidents reports) to avoid showing inappropriate ads
 - Information extraction: discard information from subjective sentences
 - Summarization: present summaries with diverse points of view

Intuitive approach: list of strongly opinionated terms

Pang et al., 2002 [13]

- Idea: manually create a list of positive and negative terms and classify documents according to how frequent positive and negative terms occur
- Assessor 1 (for movie reviews)
 - positive: dazzling, brilliant, phenomenal, excellent, fantastic
 - negative: suck, terrible, awful, unwatchatble, hideous
 - Accuracy~58%, ties~75%
- Posthoc analysis
 - Positive: love, wonderful, best, great, superb, still beautiful
 - Negative: bad, worst, stupid, waste, boring, ?, !
 - Accuracy~69%, ties~16%

gold standard false true true false true algorithm positive positive (TP) (FP) alse false true negative negative (FN) (TN)

 $accuracy = \frac{TP + TN}{TP + TN + FP + FN}$

if the test data contains the same number of true/false items, a random "classifier" on av. achieves 50% accuracy

Machine learning for opinion mining [12]

Pang et al., 2002 [13]

• Predefined set of features $\{f_1, f_2, ..., f_m\}$ that are related to the classification task at hand

number of times

- Unigram: *still*
- Bigram: really stinks
- POS tags
- #exclamation marks
- #question marks
- #🙂
- #🕱

•

$$D = \begin{pmatrix} n_1(D) \\ n_2(D) \\ n_3(D) \\ \dots \\ n_m(D) \end{pmatrix}$$

document as a vector
m dimensions=#features

Finding good features is hard!

Naïve Bayes

- A review is assigned the class $c' = \arg \max_{c} P(c \mid D)$
- Bayes' rule

$$P(c \mid D) = \frac{P(c)P(D \mid c)}{P(D)}$$

 Features are assumed to be conditionally independent given D's class

Class prior: used when one
class is known to occur more
frequently than the other
(e.g. spam emails vs. legitimite
emails)
$$P(c \mid D) = \frac{P(c \mid D) = P(c \mid C)}{P(D)}$$
 add-one smoothing
(to avoid zero prob.)

Results

Pang et al., 2002 [13]

- Dataset: 700 positive and 700 negative reviews from IMDB
- Negation tagging: addition of NOT_ to every word between a negation indicator and a punctuation mark
- 3-fold cross validation

	#features	freq. or pres.	accuracy
unigrams	16165	freq.	78.7
unigrams	16165	pres.	81.0
Unigrams+bigrams	32330	pres.	80.6
Unigrams+POS	16695	pres.	81.5
adjectives	2633	pres.	77.0

Sentiment anlysis challenges

Sentiment polarity

- The movie is not bad.
- Modeling sequential information
 - A is better than B. vs. B is better than A.
- Subtle sentiments
 - The polar express seems overly concerned with aping real life instead of creating its own universe.
- Distinction between opinions and facts
 - I love this movie. Vs. This is a love story.
- Documents often contain positive and negative opinions
 - The camera in the phone is great, but the radio is rubbish.
- Even more: cross-language mining, microblogs, emotions

Spam

Spam detection

• Spam: unsolicited (and possibly commercial) bulk messages

• Types of spam

- Email
- Instant messaging (spim)
- Internet telephony (spit)
- Mobile phone
- Web

Motivation: monetary

Web spam

 Spammers cannot send pages directly to the user, need "cooperation" of search engines (spamdexing)

- Users tend to click only on the top ranked results of search engine listings, thus spammers need to create pages that score highly
- Grey area between ethical Search Engine Optimization (SEO) and unethical spam
- Adversarial relationship between web site administrator and search engine administrator

Formally

Spamming [14]

- "Any deliberate human action that is meant to trigger an unjustifiably favorable relevance or importance for some web page, considering the page's true value"
- Spamming [15]
 - "A web page created for the sole purpose of attracting search engine referrals (to this page or some other target page)."

Web spam classification

• Content spam: making changes to a web page's content [14, 9]

- Dumping
- Weaving
- Phrase stitching
- Keyword stuffing
- ...
- Topological spam: spamming with the help of link farms
- Cloaking [16]: "... a hiding technique used by some Web servers to deliver one page to a search engine for indexing while serving an entirely different page to users browsing the site."

Content spam

- Can be automatically detected with high accuracy [9]
- Features

TUDelft

- Number of words in the page
- Number of words in the page title
- Average lenth of the words
- Amount of anchor text
- Fraction of the N most popular words in the page
- Fraction of visible content
- Compression ratio

num. words/title

Source: [9]

Decision trees

A detour

- Widely used machine learning approach
- Learned function represented in a decision tree
- Human readable
- Robust to noise

TUDelft

• Can also be represented as if/then rules

PlayTennis={YES,NO}

Decision tree overview taken from: Machine Learning, Tom Mitchell, McGraw Hill, 1997.

Decision trees

A detour

- Widely used machine learning approach
- Learned function represented in a decision tree
- Human readable
- Robust to noise
- Can also be represented as if/then rules

PlayTennis={YES,NO}

 $I = \{outlook_{rain}, wind_{strong}, humidity_{high}\}$

Decision trees: appropriate when ...

A detour

- Instances are represented by attribute-value pairs
 - Either discrete (e.g. outlook={sunny,overcast,rain}) or real-valued
- Target function has discrete output values (e.g. *PlayTennis={YES,NO}*)
- The training data may contain errors
- The training data may contain missing attribute values

Decision trees

A detour

No back-tracking

• General idea: top-down greedy built-up of the tree

① Which attribute should be the tree's root/new node?

- Teach each attribute alone, how well does it classify the training examples? Pick the best one.
- 2 Child node created for each possible value of this attribute
- ③ Sort training instances accordingly
- ④ Go to step (1), repeat until all training instances classified

Decision trees: what is the 'best' attribute?

A detour

•
$$Entropy(S) = \sum_{i=1}^{r} -p_i \log_2 p_i$$

• Example: 2 classes, 14 instances, 9+ and 5-

$$Entropy(S) = -p_{+} \log_2 p_{+} - p_{-} \log_2 p_{-}$$

$$= -\frac{9}{14}\log_2\frac{9}{14} - \frac{5}{14}\log_2\frac{5}{14} = 0.94$$

• **Information gain**: expected reduction in entropy when partitioning the training data according to the attribute

$$Gain(S,A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

define: $0 \log 0 = 0$

Decision trees: what is the 'best' attribute?

A detour

Decision trees: there is more to it

A detour

- Algorithm: ID3
 - Improvements exist
 - Pruning, etc.
- Most famous decision tree algorithm: C4.5
- Continuous variables can be included via boolean attributes
 - *A_c* is true if *A*<*c*, false otherwise
- Tom Mitchell: http://www.cs.cmu.edu/~tom/mlbook.html

Classification accuracy

Ntoulas et al., 2006 [9]

Detection rates

Topological spam

Becchetti et al., 2006 [10]

 Link farm: a densely connected set of pages, created explicitly with the purpose of deceiving a link-based ranking algorithm

Becchetti et al., 2006 [10]

• Distribution of in-degree and out-degree

• δ measures the maximum difference of the cumulative distribution functions; the larger the more different the distributions

Becchetti et al., 2006 [10]

- Degree/degree ratio: ratio between a node's degree and the average degree of its neighbours (in-links and out-links)
- PageRank

TUDelft

47

Becchetti et al., 2006 [10]

- TrustRank: given a seed of trusted nodes, propagate their scores along the edges
- Intuition: a page with high PageRank but without a relationship with a trusted page is suspicious
- Trusted nodes: pages listed in the Open Directory Project (manually maintained)
- Estimated non-spam mass: amount of PageRank score that a page receives from trusted pages

Becchetti et al., 2006 [10]

- 163 attributes in total
- Decision tree classifier
- Detection rate: 80%

spam sites classified as spam

spam sites

• False positive rate: 1-3%

#normal sites classified as spam

#normal sites

Ongoing issues in spam detection

- Spam detection is an arms race between the spammer and the search engine
- Unsupervised spam filtering
 - Shown approaches require training data (ham&spam), which is costly and ineffective (spammers learn)
- Unsupervised approaches do not rely on training data
- Image spam (most spam detection methods based on textual content)

Folksonomies

Folksonomies

Delft stadhuis Markt

This photo was taken on May 24, 2011.	
Illini 71 views 4.4 10 comments	
This photo belongs to	
Gerard Stolk (vers le Midi Carême)'s photos (7,404)	
This photo also appears in	
This photo also appears in voorjaar 2011 (set)	
This photo also appears in voorjaar 2011 (set) Tags	Delft stadhuis
This photo also appears in voorjaar 2011 (set) Tags Delft • stadhuls • Markt	> Delft stadhuis
This photo also appears in voorjaar 2011 (set) Tags Delft • stadhuls • Markt	> Delft stadhuis Markt
This photo also appears in voorjaar 2011 (set) Tags Delft • stadhuis • Markt License	> Delft stadhuis Markt

Folksonomies

Social tagging systems

- What: services allow the labeling of content with keywords (tags) chosen freely by the users
 - **Narrow** tagging rights: tagging resources limited to one/several users (e.g. Flickr)
 - **Broad** tagging rights: tagging of resources by the entire community (e.g. Bibsonomy)
- Why: content collection, content management and sharing
- Alternative to existing top-down categorization techniques (taxonomy, dictionary)
- Folksonomy: social tagging data, social tagging communities

Advantages & disadvantages

Advantages

- Easily accessible to users (adding tags is simple)
- Stay up-to-date as long as user provides tags
- Items can belong to many "categories"
- Self-moderation
- Disadvantages
 - Flat structure
 - Uncontrolled vocabulary which leads to ambiguity
 - Lack of precision
 - Noise (spam, malicious users)

Tag ambiguity on Flickr

Searching for "apple" on Flickr (14.03.2012)

Sort: Relevant | Recent | Interesting

From Mikmac

From BrownSuga'

From je@n

RENALIN

From Kanko*

From { karen }

From SKT ...

From These ...

From Diica

From linastyle

From je@n

From sophist1c...

From The Pug...

From frado76

From flyzor

From Goran Aničić

From je@n

From emarschn

View: Small | Medium | Detail | Slideshow

Applications of Folksonomies

• Item recommendation

- Recommend new interesting items to folksonomy users based on their profile
- Tag recommendation
 - Aid user who adds a resource to the system by suggesting potentially matching tags
- Personalized web search

Formally

Vallet et al., 2010 [11]

- Folksonomy *F* can be defined as a tuple *F*={*T*,*U*,*D*,*A*}
- $T = \{t_1, \dots, t_L\}$ is the set of tags expressed by the folksonomy
- $U = \{u_1, \dots, u_M\}$ is the set of users that annotate documents
- $D = \{d_1, ..., d_n\}$ is the set of documents that are annotated with tags T
- The set of annotations of each tag t_l to document d_n by user u_m

$$A = \{(u_m, t_l, d_n)\} \in U \times T \times D$$

User and document profiles

Vallet et al., 2010 [11]

• User profile

 $\vec{u}_m = (u_{m,1}, ..., u_{m,L}), \text{ where } u_{m,l} = |\{(u_m, t_l, d) \in A \mid d \in D\}|$ Based on the number of times the user has tagged documents with tag t_i

• Document profile

 $\vec{d}_n = (d_{n,1}, \dots, d_{n,L}), \text{ where } d_{n,l} = |\{(u, t_l, d_n) \in A \mid u \in U\}|$

Based on the number of times the document has been tagged with tag t_i

Personalization step

Vallet et al., 2010 [11]

• Given query *Q*, the top *s* documents retrieved by a search system are re-ranked according to the user and document profiles

$$tf if(u_m, d_n) = \sum_{l} \left(tf_{u_m}(t_l) \times iuf(t_l) \times tf_{d_n}(t_l) \times idf(t_l) \right)$$

user-baseduser-based inversedocument-baseddocument-basedtag frequencytag frequencytag frequencyinverse tag frequency

Evaluation

Vallet et al., 2010 [11]

- Document *d* is relevant to user *u* if it occurs in his profile
- Tagging information of each user is split into user profile partition and automatic topic generation partition

Evaluation

Vallet et al., 2010 [11]

- Topics are generated from each document in the test partition
 - ① Extract the k most popular tags of d
 - 2 Use extracted tags as query to a Web search engine and return the top *R* results
 - \bigcirc If *d* is not in *R*, discard the topic
 - 4 Apply personalization to the result list & report reciprocal rank of d 1

$$RR = \frac{1}{r}$$

Results

Vallet et al., 2010 [11]

- Social tagging sytem: Delicious
- 2000 users: 100 bookmarks per user, 90% for profile creation and 10% for topic generation
- 161542 documents and 69930 distinct tags
- Web search engine: Yahoo!

	Web search	Web search + tf.if
MRR	0.329	0.402
P@5	0.452	0.565
P@10	0.579	0.690
P@20	0.708	0.798

What else can be done with folksonomies?

- Rank photos based on their attractiveness by exploiting the community feedback on Flickr; learn a classifier based on assigned tags (ugly, beautiful, gross, awesome, ...)
- Locate expert users in social tagging systems based on the quality of his resources
- Predict the latitude/longitude of a photo based on Flickr tags

Summary

- Web structure
- Evolution of the Web
- Personalization on the Web
- Fighting the spammers!

Sources

- 1) Introduction to Information Retrieval. Manning et al. 2008.
- 2) Information retrieval. Keith van Rijsbergen, 1979.
- 3) Managing gigabytes, Witten et al. 1999.
- 4) Graph structure in the Web. Broder et al. 1999.
- 5) The PageRank Citation Ranking: Bringing Order to the Web. Page et al. 1999.
- 6) A taxonomy of Web search. Broder. 2002.
- 7) The indexable Web is more than 11.5 billion pages. Gulli & Signorini. 2005.
- 8) A Large-Scale Study of the Evolution of Web Pages. Fetterly et al. 2003.
- 9) Detecting spam web pages through content analysis. Ntoulas et al. 2006.
- 10) Link-based characterization and detection of web spam. Becchetti et al. 2006.
- 11) Personalizing web search with folksonomy-based user and document profiles. Vallet et al. 2010.
- 12) Authoritative Sources in a Hyperlinked Environment. Kleinberg. 1998.
- 13) Thumbs up? Sentiment classification using machine learning techniques. Pang et al. 2002.
- 14) Web spam taxonomy. Gyöngyi et al. 2004
- 15) Spam, damn spam, and statistics: using statistical analysis to locate spam web pages. Fetterly et al. 2004.
- 16) Improving cloaking detection using search query popularity and monetizability. Chellapilla et al. 2006

